

Justification of Logarithmic Loss via the Benefit of Side Information

Jiantao Jiao, *Student Member, IEEE*, Thomas Courtade, *Member, IEEE*, Kartik Venkat, *Student Member, IEEE*, and Tsachy Weissman, *Fellow, IEEE*

Abstract

We consider a natural measure of the benefit of side information: the reduction in optimal estimation risk when side information is available to the estimator. When such a measure satisfies a natural data processing property, and the source alphabet has cardinality greater than two, we show that it is uniquely characterized by the optimal estimation risk under logarithmic loss, and the corresponding measure is equal to mutual information. Further, when the source alphabet is binary, we characterize the only admissible forms the measure of predictive benefit can assume. These results allow to unify many of the causality measures in the literature as instantiations of directed information, and present the first axiomatic characterization of mutual information without requiring the sum or recursivity property.

I. INTRODUCTION

In statistical decision theory, it is often a controversial issue to choose the appropriate loss function in quantifying the risk for a given application. One popular loss function is called *logarithmic loss*, defined as follows. Let \mathcal{X} be a finite set with $|\mathcal{X}| = n$, let Γ_n denote the set of probability measures on \mathcal{X} , and let $\bar{\mathbb{R}}$ denote the extended real line.

Definition 1 (Logarithmic Loss). *Logarithmic loss $\ell_{\log}: \mathcal{X} \times \Gamma_n \rightarrow \bar{\mathbb{R}}$ is defined by*

$$\ell_{\log}(x, P) = -\log P(x), \quad (1)$$

Jiantao Jiao, Kartik Venkat, and Tsachy Weissman are with the Department of Electrical Engineering, Stanford University.
Email: {jiantao,kvenkat,tsachy}@stanford.edu

Thomas Courtade is with the Department of Electrical Engineering and Computer Sciences, University of California, Berkeley.
Email: courtade@eecs.berkeley.edu

This work was supported in part by the NSF Center for Science of Information under grant agreement CCF-0939370.

where $P(x)$ denotes the probability of x under measure P .

Logarithmic loss has enjoyed numerous applications in various fields. For instance, its usage in statistics dates back to Good [1], and it has found a prominent role in learning and prediction (cf. Cesa-Bianchi and Lugosi [2, Ch. 9]). Logarithmic loss also assumes an important role in information theory, where many of the fundamental quantities (e.g., entropy, relative entropy, etc.) can be interpreted as the optimal estimation risk under logarithmic loss. The use of the logarithm in defining entropy arises due to its various axiomatic characterizations, the first of which dates back to Shannon [3].

The main contribution of this paper is in providing fundamental justification for inference using logarithmic loss. In particular, we show that a single modest and natural Data Processing requirement mandates the use of logarithmic loss. We begin by posing the following.

Question 1 (Benefit of Side Information). *Suppose X, Z are jointly distributed random variables. How significant is the contribution of Z for inference on X ?*

II. PROBLEM FORMULATION AND MAIN RESULTS

Toward answering Question 1, let $\ell: \mathcal{X} \times \mathcal{Y} \rightarrow \bar{\mathbb{R}}$ be an arbitrary loss function with reproduction alphabet \mathcal{Y} , where \mathcal{Y} is arbitrary. Given $(X, Z) \sim P_{XZ}$, it's natural to quantify the benefit of additional side information Z by computing the difference between the expected losses in estimating $X \in \mathcal{X}$ with and without side information Z , respectively. This motivates the definition:

$$C(\ell, P_{XZ}) \triangleq \inf_{y_1 \in \mathcal{Y}} \mathbb{E}_P[\ell(X, y_1)] - \inf_{Y_2(Z)} \mathbb{E}_P[\ell(X, Y_2)], \quad (2)$$

where $y_1 \in \mathcal{Y}$ is deterministic, and $Y_2 = Y_2(Z) \in \mathcal{Y}$ is any measurable function of Z . We require that indeterminate forms like $\infty - \infty$ do not appear in the definition of $C(\ell, P_{XZ})$. By taking Z to be independent of X , we obtain for all $P \in \Gamma_n$, $|\inf_{y_1 \in \mathcal{Y}} \mathbb{E}_P[\ell(X, y_1)]| < \infty$.

The formulation (2) has appeared previously in the statistics literature. In [4], Dawid defined the *coherent dependence function*, which is equivalent to (2), and used it to quantify the dependence between two random variables X, Z . Our framework of quantifying the predictive benefit of side information is also closely connected to the notion of proper scoring rules and the literature on probability forecasting in statistics. The survey by Gneiting and Raftery [5] provides a good overview.

Having introduced the yardstick in (2), we can now reformulate the question of interest: Which loss function(s) ℓ can be used to define $C(\ell, P_{XZ})$ in a meaningful way? Of course, “meaningful” is open to interpretation, but it is desirable that $C(\ell, P_{XZ})$ be well-defined, at minimum. This motivates the following axiom:

Data Processing Axiom. For all distributions P_{XZ} , the quantity $C(\ell, P_{XZ})$ satisfies

$$C(\ell, P_{TZ}) \leq C(\ell, P_{XZ})$$

whenever $T(X) \in \mathcal{X}$ is a statistically sufficient transform of X for Z .

We remind the reader that the statement ‘ T is a statistically sufficient transform of X for Z ’ means that the following two Markov chains hold:

$$T - X - Z, \quad X - T - Z \quad (3)$$

In other words, $T(X)$ is a lossless representation of all of the information X contains about Z .

In words, the Data Processing Axiom stipulates that processing the data $X \rightarrow T$ cannot boost the predictive benefit of the side information¹.

To convince the reader that the Data Processing Axiom is a natural requirement, suppose instead that the Data Processing Axiom did not hold. Since X and T are mutually sufficient statistics for Z , this would imply that there is *no* unique value which quantifies the benefit of side information Z for the random variable of interest. Thus, the Data Processing Axiom is needed for the benefit of side information to be well-defined.

Benign as the Data Processing Axiom, it has far-reaching implications for the form $C(\ell, P_{XZ})$ can take. This is captured by our first main result:

Theorem 1. Let $n \geq 3$. Under the Data Processing Axiom, the function $C(\ell, P_{XZ})$ is uniquely determined by the mutual information,

$$C(\ell, P_{XZ}) = I(X; Z), \quad (4)$$

up to a multiplicative factor.

The following corollary immediately follows from Theorem 1.

Corollary 1. Let $n \geq 3$. Under the Data Processing Axiom, the benefit of additional side information Z for inference on X with common side information W , i.e.

$$\inf_{Y_1(W)} \mathbb{E}_P[\ell(X, Y_1)] - \inf_{Y_2(Z, W)} \mathbb{E}_P[\ell(X, Y_2)], \quad (5)$$

¹In fact, the Data Processing Axiom is weaker than this general data processing statement since it only addresses statistically sufficient transformations of X .

is uniquely determined by the conditional mutual information,

$$I(X; Z|W), \quad (6)$$

up to a multiplicative factor.

Thus, up to a multiplicative factor, we see that logarithmic loss generates the *only* measure of predictive benefit (defined according to (2)) which satisfies the Data Processing Axiom. In other words, Theorem 1 provides a definitive answer to Question 1 under the framework we have described, and also highlights the special role that logarithmic loss plays.

Theorem 1 shows that mutual information is natural to measure the amount of reduction of statistical risk when we have side information. Incidentally, Erkip and Cover [6] argued that mutual information was a natural quantity in the context of portfolio theory, where it emerges as the increase in growth rate due to the presence of side information.

It is worth mentioning that Theorem 1 is closely connected to existing results on axiomatic characterizations of information measures; see Csiszár [7] for a survey. To emphasize our contribution, we note that Csiszár [7] names only the axiomatic result of Ac zel, Forte, and Ng [8] as a characterization of information measures that requires neither recursivity nor the sum property. However, [8] focuses on entropy characterization, and the framework therein does not extend to the problem we consider.

Interestingly, the assumption that $n \geq 3$ in Theorem 1 is essential. The class of solutions for the binary alphabet setting is characterized by the following theorem.

Theorem 2. *Let $n = 2$. $C(\ell, P_{XZ})$ is of the form*

$$C(\ell, P_{XZ}) = \mathbb{E}_Z[G(P_{X|Z})] - G(P_X)$$

for a symmetric convex function $G((p, 1-p)) : \Gamma_2 \rightarrow \mathbb{R}$ if, and only if, the Data Processing Axiom holds.

It is worth mentioning that there is an interesting regime of observations surrounding the characterization of information measures, which is sensitive to the alphabet size being binary or larger. This phenomenon is explored in details in [9].

The rest of this paper is organized as follows. In Section III, we explore the connections between our results and the existing literature on causal analysis, including Granger and Sims causality, Geweke's measure, transfer entropy, and directed information. The proofs of Theorems 1 and 2 are provided in Section IV. Proofs of some auxiliary lemmas are deferred to the appendix.

III. CAUSALITY MEASURES: AN AXIOMATIC VIEWPOINT

Inferring causal relationships from observed data plays an indispensable part in scientific discovery. Granger, in his seminal work [10], proposed a predictive test for inferring causal relationships. To state his test, let X_t, Y_t, U_t be stochastic processes, where X_t, Y_t are the processes of interest, and U_t contains all information in the universe accumulated up to time t . Granger's causality test asserts that Y_t causes X_t , denoted by $Y_t \Rightarrow X_t$, if we are better able to predict X_t using the past information of U_t , than by using all past information in U_t apart from Y_t . In Granger's definition, the quality of prediction is measured by the squared error risk achieved by the optimal unbiased least-squares predictor.

In his 1980 paper, Granger [11] introduced a set of operational definitions which made it possible to derive practical testing procedures. For example, he assumes that we must be able to specify U_t in order to perform causality tests, which is slightly different from his original definition which required knowledge of all information in the universe (which is usually unavailable).

Later, Sims [12] introduced a related concept of causality, which was proved to be equivalent to Granger's definition in Sims [12], Hosoya [13], and Chamberlain [14] in a variety of settings.

Motivated by Granger's framework for testing causality using linear prediction, Geweke [15][16] proposed a causality measure to quantify the extent to which Y is causing X . Quoting Geweke (emphasis ours):

“The empirical literature abounds with tests of independence and unidirectional causality for various pairs of time series, but there have been virtually no investigations of the degree of dependence or the extent of various kinds of feedback. The latter approach is more realistic in the typical case in which the hypothesis of independence of unidirectional causality is not literally entertained, but it requires that one be able to measure linear dependence and feedback.”

In other words, Geweke makes the important distinction between a *causality test* which makes a binary decision on whether one process causes another, and a *causality measure* which quantifies the degree to which one process causes another. Geweke proposed the following measure as a natural starting point:

$$F_{Y \Rightarrow X} \triangleq \ln \frac{\sigma^2(X_t | X^{t-1})}{\sigma^2(X_t | X^{t-1}, Y^{t-1})}, \quad (7)$$

where $\sigma^2(X_t | X^{t-1}, Y^{t-1})$ is the variance of the prediction residue when predicting X_t via the optimal linear predictor constructed from observation X^{t-1}, Y^{t-1} . Note that if $F_{Y \Rightarrow X} > 0$, we could conclude $Y_t \Rightarrow X_t$ according to Granger's test.

It has long been observed that the restriction to optimal linear predictors in testing causality is not necessary. In fact, Chamberlain [14] proved a general equivalence between Granger and Sims' causality

tests by replacing linear predictors with conditional independence tests. However, the most natural generalization of (7) wasn't clear until Gourieroux, Monfort, and Renault [17] proposed the so-called *Kullback causality measures* in 1987. It is now well-known that Kullback causality measures are equivalent to (7) under linear Gaussian models (cf. Barnett, Barrett and Seth [18]).

Using information theoretic terms, Kullback causality measures are nothing but the directed information introduced by Massey [19], and motivated by Marko [20]. Using modern notation, the directed information from X^n to Y^n is defined as

$$I(X^n \rightarrow Y^n) \triangleq \sum_{i=1}^n I(X^i; Y_i | Y^{i-1}) \quad (8)$$

$$= H(Y^n) - H(Y^n \| X^n), \quad (9)$$

where $H(Y^n \| X^n)$ is the *causally conditional entropy*, defined by

$$H(Y^n \| X^n) \triangleq \sum_{i=1}^n H(Y_i | Y^{i-1}, X^i). \quad (10)$$

Massey and Massey [21] established the pleasing conservation law of directed information:

$$I(X^n; Y^n) = I(X^n \rightarrow Y^n) + I(Y^{n-1} \rightarrow X^n) \quad (11)$$

$$\begin{aligned} &= I(X^{n-1} \rightarrow Y^n) + I(Y^{n-1} \rightarrow X^n) \\ &\quad + \sum_{i=1}^n I(X_i; Y_i | X^{i-1}, Y^{i-1}), \end{aligned} \quad (12)$$

which implies that the extent to which process X_t influences process Y_t and vice-versa always sum to the total mutual information between the two processes. Since $I(Y^{n-1} \rightarrow X^n)$ can be expressed as

$$I(Y^{n-1} \rightarrow X^n) = \sum_{i=1}^n H(X_i | X^{i-1}) - H(X_i | X^{i-1}, Y^{i-1}),$$

X_i being conditionally independent of Y^{i-1} given X^{i-1} is equivalent to $I(Y^{n-1} \rightarrow X^n) = 0$. This corresponds precisely to the definition of general Granger non-causality. Permuter, Kim, and Weissman [22] showed various applications of directed information in portfolio theory, data compression, and hypothesis testing in the presence of causality constraints.

We remark that, for practical applications, the directed information between stochastic processes can be computed using the universal estimators proposed in [23], which exhibit optimal statistical and convergence properties.

Finally, we note that the notion of *transfer entropy* in the physics literature, which was proposed by Schreiber [24] in 2000, turns out to be equivalent to directed information.

To connect our present discussion on causality measures to Theorem 1, we recall that the directed information rate [25] between a pair of jointly stationary finite-alphabet processes X_t, Y_t can be written as:

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{1}{n} I(Y^{n-1} \rightarrow X^n) \\ = \inf_{T_1(X_{-\infty}^{-1})} \mathbb{E}[\ell_{\log}(X_0, T_1)] - \inf_{T_2(X_{-\infty}^{-1}, Y_{-\infty}^{-1})} \mathbb{E}[\ell_{\log}(X_0, T_2)]. \end{aligned}$$

In light of this, we can conclude from Theorem 1 and Corollary 1 that the directed information rate is the *unique* measure of causality which assumes the form (2) and satisfies the Data Processing Axiom. Thus, our axiomatic viewpoint explains why the same causality measure has appeared so often in varied fields including economics, statistics, information theory, and physics. Except in the binary case, we roughly have the following: *All reasonable causality measures defined by a difference of predictive risks must coincide.*²

IV. PROOF OF MAIN RESULTS

In this section, we provide complete proofs of Theorems 1 and 2 and highlight the key ideas.

To begin, we need to put all estimation risk on the common footing by eliminating the arbitrary nature of the reconstruction alphabet \mathcal{Y} . The following lemma achieves this goal.

Lemma 1. *There exists a bounded convex function $V : \Gamma_n \rightarrow \mathbb{R}$, depending on ℓ , such that*

$$C(\ell, P_{XZ}) = \left(\sum_z P_Z(z) V(P_{X|Z}) \right) - V(P_X). \quad (13)$$

The proof of Lemma 1 follows from defining $V(P)$ by

$$V(P) = - \inf_{y \in \mathcal{Y}} \mathbb{E}_P[\ell(X, y)], \quad (14)$$

and its details are deferred to the appendix. In the statistics literature, $V(P)$ is usually called the *negative generalized entropy*, and we refer to Dawid [26] for details.

In the literature of concentration inequalities, the following functional

$$H_\Phi(Z) = \mathbb{E}\Phi(Z) - \Phi(\mathbb{E}Z), \quad (15)$$

²Here, the authors' interpretation of "reasonable" is reflected by the Data Processing Axiom. In the context of this section, the Data Processing Axiom stipulates that any reasonable causality measure should be invariant under statistically sufficient transformations of the data – a desirable property and natural criterion.

where Φ is a convex function, is called Φ -entropy. As shown in Lemma 1, functional $C(\ell, P_{XZ})$ is closely related to the notion of Φ -entropy. We refer to Boucheron, Lugosi, and Massart [27, Ch. 14] for a nice survey on the usage of Φ -entropies in proving concentration inequalities.

The next lemma asserts that we only need to consider symmetric $V(P)$.

Lemma 2. *Under the Data Processing Axiom, there exists a symmetric finite convex function $G : \Gamma_n \rightarrow \mathbb{R}$, such that*

$$C(\ell, P_{XZ}) = \left(\sum_z P_Z(z) G(P_{X|Z}) \right) - G(P_X), \quad (16)$$

and $G(\cdot)$ is equal to $V(\cdot)$ in Lemma 1 up to a linear translation:

$$G(P) = V(P) + \langle c, P \rangle, \quad (17)$$

where $c \in \mathbb{R}^n$ is a constant vector.

The proof of Lemma 2 follows by applying a permutation to the space \mathcal{X} and applying the Data Processing Axiom. Details are deferred to the appendix.

Now we are in a position to begin the proof of Theorem 1 in earnest.

A. The case $n \geq 3$

Let $Z \in \{1, 2\}$, with $\alpha \triangleq \mathbb{P}\{Z = 1\}$. Take $P_{\lambda_1}^{(t)}, P_{\lambda_2}^{(t)}$ to be two probability vectors parametrized in the following way:

$$P_{\lambda_1}^{(t)} = (\lambda_1 t, \lambda_1(1-t), r - \lambda_1, p_4, \dots, p_n) \quad (18)$$

$$P_{\lambda_2}^{(t)} = (\lambda_2 t, \lambda_2(1-t), r - \lambda_2, p_4, \dots, p_n), \quad (19)$$

where $r \triangleq 1 - \sum_{i \geq 4} p_i$, $t \in [0, 1]$, $\lambda_1 < \lambda_2$.

Taking $P_{X|1} \triangleq P_{\lambda_1}^{(t)}$, $P_{X|2} \triangleq P_{\lambda_2}^{(t)}$, we have

$$C(\ell, P_{XZ}) = \alpha V(P_{\lambda_1}^{(t)}) + (1 - \alpha) V(P_{\lambda_2}^{(t)}) - V(\alpha P_{\lambda_1}^{(t)} + (1 - \alpha) P_{\lambda_2}^{(t)}). \quad (20)$$

Note that for any $\alpha, t, \lambda_1, \lambda_2$, the following transformation is sufficient for Z .

$$T(X) = \begin{cases} 1 & X \in \{x_1, x_2\} \\ X & \text{otherwise} \end{cases} \quad (21)$$

The Data Processing Axiom implies that for all $\alpha \in [0, 1]$ and legitimate $\lambda_2 > \lambda_1 \geq 0$,

$$\begin{aligned} & \alpha V(P_{\lambda_1}^{(t)}) + (1 - \alpha)V(P_{\lambda_2}^{(t)}) - V(\alpha P_{\lambda_1}^{(t)} + (1 - \alpha)P_{\lambda_2}^{(t)}) \\ &= \alpha V(P_{\lambda_1}^{(1)}) + (1 - \alpha)V(P_{\lambda_2}^{(1)}) - V(\alpha P_{\lambda_1}^{(1)} + (1 - \alpha)P_{\lambda_2}^{(1)}). \end{aligned} \quad (22)$$

Fixing p_4, p_5, \dots, p_n , we define the function

$$R(\lambda, t; p_4, p_5, \dots, p_n) \triangleq V(P_{\lambda}^{(t)}), \quad (23)$$

where, for notational simplicity, we denote $R(\lambda, t; p_4, p_5, \dots, p_n)$ by $R(\lambda, t)$.

Note that by definition,

$$R(\alpha\lambda_1 + (1 - \alpha)\lambda_2, t) = V(\alpha P_{\lambda_1}^{(t)} + (1 - \alpha)P_{\lambda_2}^{(t)}), \quad (24)$$

hence we know that

$$\begin{aligned} & \alpha R(\lambda_1, t) + (1 - \alpha)R(\lambda_2, t) - R(\alpha\lambda_1 + (1 - \alpha)\lambda_2, t) \\ &= \alpha R(\lambda_1, 1) + (1 - \alpha)R(\lambda_2, 1) - R(\alpha\lambda_1 + (1 - \alpha)\lambda_2, 1). \end{aligned} \quad (25)$$

Note that if we define $\tilde{R}(\lambda, t) \triangleq R(\lambda, t) - \lambda U(t) - F(t)$, where $U(t), F(t)$ are arbitrary real-valued functions, for all λ_1, λ_2, t we have

$$\begin{aligned} & \alpha R(\lambda_1, t) + (1 - \alpha)R(\lambda_2, t) - R(\alpha\lambda_1 + (1 - \alpha)\lambda_2, t) \\ &= \alpha \tilde{R}(\lambda_1, t) + (1 - \alpha)\tilde{R}(\lambda_2, t) - \tilde{R}(\alpha\lambda_1 + (1 - \alpha)\lambda_2, t), \end{aligned} \quad (26)$$

which implies that

$$\begin{aligned} & \alpha \tilde{R}(\lambda_1, t) + (1 - \alpha)\tilde{R}(\lambda_2, t) - \tilde{R}(\alpha\lambda_1 + (1 - \alpha)\lambda_2, t) \\ &= \alpha \tilde{R}(\lambda_1, 1) + (1 - \alpha)\tilde{R}(\lambda_2, 1) - \tilde{R}(\alpha\lambda_1 + (1 - \alpha)\lambda_2, 1). \end{aligned} \quad (27)$$

Taking $\lambda_1 = 0, \lambda_2 = r = 1 - \sum_{i \geq 4} p_i$, we can choose the functions $U(t), F(t)$ in a way such that

$$\tilde{R}(0, t) = A(p_4, \dots, p_n), \quad \tilde{R}(r, t) = B(p_4, \dots, p_n), \quad \forall t \in [0, 1], \quad (28)$$

where A, B are some functions of (p_4, \dots, p_n) .

Plugging (28) into (27), we know that

$$\begin{aligned} & \alpha A(p_4, \dots, p_n) + (1 - \alpha)B(p_4, \dots, p_n) - \tilde{R}((1 - \alpha)r, t) \\ &= \alpha A(p_4, \dots, p_n) + (1 - \alpha)B(p_4, \dots, p_n) - \tilde{R}((1 - \alpha)r, 1), \end{aligned}$$

which implies that

$$\tilde{R}((1 - \alpha)r, t) = \tilde{R}((1 - \alpha)r, 1), \quad \forall \alpha \in [0, 1], t \in [0, 1] \quad (29)$$

In other words, there exists a function $E : [0, 1] \rightarrow \mathbb{R}$, such that

$$\tilde{R}(\lambda, t) = E(\lambda). \quad (30)$$

Since $R(\lambda, t) = \tilde{R}(\lambda, t) + \lambda U(t) + F(t)$, we know that there exist real-valued functions E, U, F (indexed by p_4, \dots, p_n) such that

$$R(\lambda, t) = F(t) + \lambda U(t) + E(\lambda). \quad (31)$$

Expressing λ, t in terms of p_1, p_2 , we have

$$\lambda = p_1 + p_2, \quad t = \frac{p_1}{p_1 + p_2}. \quad (32)$$

By definition of $R(\lambda, t)$, we have

$$\begin{aligned} & V(p_1, p_2, p_3, p_4, \dots, p_n) \\ &= F\left(\frac{p_1}{p_1 + p_2}; p_4, \dots, p_n\right) + (p_1 + p_2)U\left(\frac{p_1}{p_1 + p_2}; p_4, \dots, p_n\right) + E(p_1 + p_2; p_4, \dots, p_n). \end{aligned} \quad (33)$$

By Lemma 2, we know that there exists a symmetric finite convex function $G : \Gamma_n \rightarrow \mathbb{R}$, such that

$$G(P) = V(P) + \langle c, P \rangle. \quad (34)$$

Now we cite a result by Gale, Klee, and Rockafellar on properties of bounded convex functions on polytopes.

Lemma ([28, Gale-Klee-Rockafellar's Theorem]). *If D is boundedly polyhedral and ϕ is a convex function on the relative interior of D which is bounded on bounded sets, then ϕ can be extended in a unique way to a continuous convex function on D .*

Taking $D = \Gamma_n$, $\phi = G$, it follows from Lemma 1 and Lemma 2 that ϕ is bounded. Gale-Klee-Rockafellar's Theorem implies that $G(P)$ can be extended in a unique way to a continuous convex function on Γ_n .

Taking $p_1 = xa$, $p_2 = x(1 - a)$, $a \in [0, 1]$, $p_3 = 1 - \left(\sum_{i \geq 4} p_i\right) - x$, and letting $x \downarrow 0$, it follows from the general expression of $G(P)$ that

$$\lim_{x \downarrow 0} G(P) = F(a; p_4, \dots, p_n) + \lim_{x \downarrow 0} E(x; p_4, \dots, p_n) + c_3(1 - \sum_{i \geq 4} p_i) + \sum_i c_i p_i. \quad (35)$$

Equation (35) implies that if F is not identically constant, $\lim_{x \downarrow 0} G(P)$ is going to depend on how we approach the boundary point $(0, 0, 1 - \sum_{i \geq 4} p_i, p_4, \dots, p_n)$, which would contradict the Gale–Klee–Rockafellar Theorem. Thus we know that $F \equiv \text{const.}$ Without loss of generality, we take $F \equiv 0$.

In other words, we have proved that G is of the form

$$G(P) = (p_1 + p_2)U\left(\frac{p_1}{p_1 + p_2}; p_4, \dots, p_n\right) + E(p_1 + p_2; p_4, \dots, p_n) + \langle c, P \rangle. \quad (36)$$

For notational simplicity, we define

$$Y(p_1, p_2; p_4, \dots, p_n) \triangleq G(P), \quad (37)$$

and denote $Y(p_1, p_2; p_4, \dots, p_n)$ by $Y(p_1, p_2)$. This gives

$$Y(p_1, p_2) = (p_1 + p_2)U\left(\frac{p_1}{p_1 + p_2}\right) + E(p_1 + p_2) + c_1 p_1 + c_2 p_2 + c_3(r - p_1 - p_2). \quad (38)$$

Since $G(P)$ is a symmetric function, we know that if we exchange p_1 and p_3 in $G(P)$, the value of $G(P)$ will not change. In other words, for $r = p_1 + p_2 + p_3$, we have

$$\begin{aligned} & (r - p_3)U\left(\frac{p_1}{r - p_3}\right) + E(r - p_3) + c_1 p_1 + c_2 p_2 + c_3 p_3 \\ &= (r - p_1)U\left(\frac{p_3}{r - p_1}\right) + E(r - p_1) + c_1 p_3 + c_2 p_2 + c_3 p_1, \end{aligned} \quad (39)$$

which is equivalent to

$$\begin{aligned} & (r - p_3)U\left(\frac{p_1}{r - p_3}\right) + E(r - p_3) + (c_3 - c_1)p_3 \\ &= (r - p_1)U\left(\frac{p_3}{r - p_1}\right) + E(r - p_1) + (c_3 - c_1)p_1. \end{aligned} \quad (40)$$

Define $\tilde{E}(x; p_4, \dots, p_n) \triangleq E(r - x; p_4, \dots, p_n) + (c_3 - c_1)x$, we have

$$(r - p_3)U\left(\frac{p_1}{r - p_3}\right) + \tilde{E}(p_3) = (r - p_1)U\left(\frac{p_3}{r - p_1}\right) + \tilde{E}(p_1). \quad (41)$$

Now we cite a result on generalizations of the so-called *fundamental equation of information theory*:

Lemma ([29][30][31]). *The most general measurable solution of*

$$f(x) + (1 - x)g\left(\frac{y}{1 - x}\right) = h(y) + (1 - y)k\left(\frac{x}{1 - y}\right), \quad (42)$$

for $x, y \in [0, 1)$ with $x + y \in [0, 1]$, where $f, h : [0, 1) \rightarrow \mathbb{R}$ and $g, k : [0, 1] \rightarrow \mathbb{R}$, has the form

$$f(x) = aH_2(x) + b_1x + d, \quad (43)$$

$$g(y) = aH_2(y) + b_2y + b_1 - b_4, \quad (44)$$

$$h(x) = aH_2(x) + b_3x + b_1 + b_2 - b_3 - b_4 + d, \quad (45)$$

$$k(y) = aH_2(y) + b_4y + b_3 - b_2, \quad (46)$$

for $x \in [0, 1], y \in [0, 1]$, where $H_2(x) = -x \ln x - (1 - x) \ln(1 - x)$ is the binary Shannon entropy and a, b_1, b_2, b_3, b_4 , and d are arbitrary constants.

Remark 1. If $f = g = h = k$ in (43)-(46), the corresponding functional equation is called the ‘fundamental equation of information theory’.

In order to apply the above lemma to our setting, we define

$$q_i = p_i/r, \quad i = 1, 2, 3 \quad (47)$$

and $h(x) = \tilde{E}(rx)/r$. Then we know

$$(1 - q_3)U\left(\frac{q_1}{1 - q_3}\right) + h(q_3) = (1 - p_1)U\left(\frac{q_3}{1 - q_1}\right) + h(q_1). \quad (48)$$

Applying the general solution of (42), setting $f = h, g = k = U$, we have

$$b_1 = b_3, b_2 = b_4. \quad (49)$$

Thus,

$$h(x) = aH_2(x) + b_1x + d, \quad (50)$$

$$U(y) = aH_2(y) + b_2y + b_1 - b_2. \quad (51)$$

By the definition of $h(x)$ and $\tilde{E}(x)$, we have that

$$E(x) = raH_2(x/r) + (b_1 + c_1 - c_3)(r - x) + d. \quad (52)$$

Plugging the general solutions to $U(x), E(x)$ into (38), and redefining the constants, we have

$$Y(p_1, p_2) = A(p_1 \ln p_1 + p_2 \ln p_2 + (r - p_1 - p_2) \ln(r - p_1 - p_2)) + Bp_1 + Cp_2 + D. \quad (53)$$

Note that the constants A, B, C, D are, in fact, functions of p_4, \dots, p_n . Therefore, we have the following general representation of the symmetric function $G(P)$:

$$\begin{aligned} G(P) &= A(p_4, \dots, p_n) (p_1 \ln p_1 + p_2 \ln p_2 + p_3 \ln p_3) \\ &\quad + B(p_4, \dots, p_n)p_1 + C(p_4, \dots, p_n)p_2 + D(p_4, \dots, p_n). \end{aligned} \quad (54)$$

Exchanging p_1, p_2 , we obtain that $B \equiv C$. Exchanging p_1, p_3 , we obtain that $B \equiv C \equiv 0$. Doing an arbitrary permutation on p_4, \dots, p_n , since p_1, p_2, p_3 enjoy two degrees of freedom, we know that $A(p_4, \dots, p_n), D(p_4, \dots, p_n)$ are symmetric functions.

Exchanging p_1, p_4 and comparing the coefficients for $p_2 \ln p_2$, we know that

$$A(p_4, p_5, \dots, p_n) = A(p_1, p_5, \dots, p_n), \quad (55)$$

since A is symmetric, and thus we can conclude that A is a constant. Now exchanging p_1, p_4 gives us

$$Ap_1 \ln p_1 - Ap_4 \ln p_4 = D(p_1, p_5, \dots, p_n) - D(p_4, p_5, \dots, p_n). \quad (56)$$

Taking partial derivatives with respect to p_1 on both sides of (56), we obtain

$$A(\ln p_1 + 1) = \frac{\partial}{\partial p_1} D(p_1, p_5, \dots, p_n). \quad (57)$$

Integrating on both sides with respect to p_1 , we know there exists a function f such that

$$D(p_1, p_5, \dots, p_n) = Ap_1 \ln p_1 + f(p_5, \dots, p_n). \quad (58)$$

Since D is symmetric, we further know that

$$D(p_4, \dots, p_n) = \sum_{i \geq 4} Ap_i \ln p_i. \quad (59)$$

To sum up, we have

$$G(P) = A \sum_{i=1}^n p_i \ln p_i. \quad (60)$$

To guarantee that $G(P)$ is convex, we need $A > 0$.

Plugging (60) into Lemma 2, the proof is complete.

B. The case $n = 2$

The ‘if’ part of Theorem 2 follows from Lemma 2. Savage’s representation of proper scoring rules [5] gives the ‘only if’ direction. In particular, the Savage representation asserts, for a convex function G , we can define a proper scoring rule $S_\ell(x, Q) : \mathcal{X} \times \Gamma_n \rightarrow \bar{\mathbb{R}}$ by

$$S_\ell(x, Q) \triangleq \langle G'(Q), Q \rangle - G(Q) - G'_x(Q), \quad (61)$$

where $G'(Q)$ denotes a sub gradient of $G(Q)$ at Q , and $G'_x(Q)$ is the component of $G'(Q)$ corresponding to $Q(x)$, and G is a convex function (see, e.g., [5] for details). By definition, a scoring rule S_ℓ is proper if the true distribution minimizes the expected loss. In other words, S_ℓ is proper if

$$P \in \inf_{Q \in \Gamma_n} \mathbb{E}_P[S_\ell(X, Q)]. \quad (62)$$

Taking the proper scoring rule S_ℓ as a loss function, and substituting into (2) defines a valid $C(\ell, P_{XZ})$.

REFERENCES

- [1] I. J. Good, "Rational decisions," *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 107–114, 1952.
- [2] N. Cesa-Bianchi and G. Lugosi, *Prediction, learning, and games*. Cambridge University Press, 2006.
- [3] C. E. Shannon, "A mathematical theory of communication," *The Bell System Technical Journal*, vol. 27, pp. 379–423, 623–656, 1948.
- [4] A. P. Dawid, "Coherent measures of discrepancy, uncertainty and dependence, with applications to bayesian predictive experimental design," Department of Statistical Science, University College London. <http://www.ucl.ac.uk/Stats/research/abs94.html>, Tech. Rep. 139, 1998.
- [5] T. Gneiting and A. E. Raftery, "Strictly proper scoring rules, prediction, and estimation," *Journal of the American Statistical Association*, vol. 102, no. 477, pp. 359–378, 2007.
- [6] E. Erkip and T. M. Cover, "The efficiency of investment information," *Information Theory, IEEE Transactions on*, vol. 44, no. 3, pp. 1026–1040, 1998.
- [7] I. Csiszár, "Axiomatic characterizations of information measures," *Entropy*, vol. 10, no. 3, pp. 261–273, 2008.
- [8] J. Aczél, B. Forte, and C. Ng, "Why the shannon and hartley entropies are 'natural'," *Advances in Applied Probability*, pp. 131–146, 1974.
- [9] J. Jiao, T. Courtade, A. No, K. Venkat, and T. Weissman, "Information measures: the curious case of the binary alphabet," *submitted to IEEE Transactions on Information Theory*.
- [10] C. W. Granger, "Investigating causal relations by econometric models and cross-spectral methods," *Econometrica: Journal of the Econometric Society*, pp. 424–438, 1969.
- [11] —, "Testing for causality: a personal viewpoint," *Journal of Economic Dynamics and control*, vol. 2, pp. 329–352, 1980.
- [12] C. A. Sims, "Money, income, and causality," *The American Economic Review*, vol. 62, no. 4, pp. 540–552, 1972.
- [13] Y. Hosoya, "On the granger condition for non-causality," *Econometrica*, vol. 45, no. 7, pp. 1735–36, 1977.
- [14] G. Chamberlain, "The general equivalence of granger and sims causality," *Econometrica: Journal of the Econometric Soc.*, pp. 569–581, 1982.
- [15] J. Geweke, "Measurement of linear dependence and feedback between multiple time series," *Journal of the American Statistical Association*, vol. 77, no. 378, pp. 304–313, 1982.
- [16] J. F. Geweke, "Measures of conditional linear dependence and feedback between time series," *Journal of the American Statistical Association*, vol. 79, no. 388, pp. 907–915, 1984.
- [17] C. Gourieroux, A. Monfort, and E. Renault, "Kullback causality measures," *Annales d'Economie et de Statistique*, pp. 369–410, 1987.
- [18] L. Barnett, A. B. Barrett, and A. K. Seth, "Granger causality and transfer entropy are equivalent for gaussian variables," *Physical review letters*, vol. 103, no. 23, p. 238701, 2009.
- [19] J. L. Massey, "Causality, feedback, and directed information," in *Proc. Int. Symp. Inf. Theory Appl.*, Honolulu, HI, Nov. 1990, pp. 303–305.
- [20] H. Marko, "The bidirectional communication theory—a generalization of information theory," *IEEE Trans. Comm.*, vol. 21, pp. 1345–1351, 1973.
- [21] J. L. Massey and P. C. Massey, "Conservation of mutual and directed information," in *Proc. IEEE Int. Symp. Inf. Theory*, 2005, pp. 157–158.
- [22] H. H. Permuter, Y.-H. Kim, and T. Weissman, "Interpretations of directed information in portfolio theory, data compression, and hypothesis testing," *Information Theory, IEEE Transactions on*, vol. 57, no. 6, pp. 3248–3259, 2011.

- [23] J. Jiao, H. Permuter, L. Zhao, Y.-H. Kim, and T. Weissman, “Universal estimation of directed information,” *Information Theory, IEEE Transactions on*, vol. 59, no. 10, pp. 6220–6242, 2013.
- [24] T. Schreiber, “Measuring information transfer,” *Physical review letters*, vol. 85, no. 2, p. 461, 2000.
- [25] G. Kramer, *Directed Information for Channels with Feedback*. Konstanz: Hartung-Gorre Verlag, 1998, Dr. sc. thchn. Dissertation, Swiss Federal Institute of Technology (ETH) Zurich.
- [26] A. P. Dawid, “The geometry of proper scoring rules,” *Annals of the Institute of Statistical Mathematics*, vol. 59, no. 1, pp. 77–93, 2007.
- [27] S. Boucheron, G. Lugosi, and P. Massart, *Concentration inequalities: A nonasymptotic theory of independence*. Oxford University Press, 2013.
- [28] D. Gale, V. Klee, and R. Rockafellar, “Convex functions on convex polytopes,” *Proceedings of the American Mathematical Society*, vol. 19, no. 4, pp. 867–873, 1968.
- [29] P. Kannappan and C. Ng, “Measurable solutions of functional equations related to information theory,” *Proceedings of the American Mathematical Society*, pp. 303–310, 1973.
- [30] G. Maksa, “Solution on the open triangle of the generalized fundamental equation of information with four unknown functions,” *Utilitas Math*, vol. 21, pp. 267–282, 1982.
- [31] J. Aczél and C. Ng, “Determination of all semisymmetric recursive information measures of multiplicative type on n positive discrete probability distributions,” *Linear algebra and its applications*, vol. 52, pp. 1–30, 1983.

APPENDIX

A. Proof of Lemma 1

Define $V : \Gamma_n \rightarrow \mathbb{R}$ by

$$V(P) = - \inf_{y \in \mathcal{Y}} \mathbb{E}_P[\ell(X, y)]. \quad (63)$$

Since $\mathbb{E}_P[\ell(X, y)]$ is linear in P , and $V(P)$ is the pointwise supremum over a family of linear functions of P , we know $V(P)$ is convex and lower semi-continuous on Γ_n .

Since Γ_n is a compact set, we know that the lower semi-continuous function $V(P)$ attains its minimum on Γ_n .

At the same time, since Γ_n is a polytope, we know $\forall P = (p_1, p_2, \dots, p_n) \in \Gamma_n$, we have $P = \sum_{i=1}^n p_i \delta_i$, where $\delta_i = (0, 0, \dots, 1, 0, \dots, 0)$ is a distribution that puts mass one at symbol i .

Since $V(P)$ is convex, we have

$$V(P) = V\left(\sum_{i=1}^n p_i \delta_i\right) \leq \sum_{i=1}^n p_i V(\delta_i) \leq \max\{V(\delta_i), 1 \leq i \leq n\}. \quad (64)$$

That is to say, the function $V(P)$ attains its maximum at one of the boundary points δ_i . Thus, we know that $V(P)$ is bounded.

Now we proceed to show that

$$\inf_{Y(Z)} \mathbb{E}_P[\ell(X, Y)] = - \sum_z P_Z(z) V(P_{X|Z}). \quad (65)$$

We first argue that $\inf_{Y(Z)} \mathbb{E}_P[\ell(X, Y)] \geq -\sum_z P_Z(z)V(P_{X|Z})$ holds, then we argue strict inequality is not possible.

By definition of the infimum, there exists a sequence of measurable functions $\{Y_n(z)\}$ such that

$$\lim_{n \rightarrow \infty} \mathbb{E}_P[\ell(X, Y_n(Z))] = \inf_{Y(Z)} \mathbb{E}_P[\ell(X, Y)]. \quad (66)$$

However, by the law of iterated expectation, we have

$$\mathbb{E}_P[\ell(X, Y_n(Z))] = \mathbb{E}_P[\mathbb{E}_P[\ell(X, Y_n(Z))|Z]] \quad (67)$$

$$\geq \mathbb{E}_P[-V(P_{X|Z})] \quad (68)$$

$$= -\sum_z P_Z(z)V(P_{X|Z}). \quad (69)$$

Letting $n \rightarrow \infty$ on both sides, we know

$$\inf_{Y(Z)} \mathbb{E}_P[\ell(X, Y)] \geq -\sum_z P_Z(z)V(P_{X|Z}). \quad (70)$$

Now suppose for some P_{XZ} it holds strict inequality in (70). By definition of $V(P_{X|Z})$, for all z , there exists a sequence $\{y_n(z)\} \subset \mathcal{Y}$ such that

$$V(P_{X|Z}(\cdot|z)) = -\inf_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} P_{X|Z}(x|z)\ell(x, y) \quad (71)$$

$$= -\lim_{n \rightarrow \infty} \sum_{x \in \mathcal{X}} P_{X|Z}(x|z)\ell(x, y_n(z)) \quad (72)$$

Without loss of generality we could assume the sequence $\sum_{x \in \mathcal{X}} P_{X|Z}(x|z)\ell(x, y_n(z))$ is non-increasing.

Now define a sequence of measurable functions of Z taking values in \mathcal{Y} as $Y_n(Z) = y_n(Z)$. We have

$$\lim_{n \rightarrow \infty} \mathbb{E}_P[\ell(X, Y_n(Z))] = \lim_{n \rightarrow \infty} \mathbb{E}_P[\mathbb{E}_P[\ell(X, Y_n(Z))|Z]] \quad (73)$$

$$= \mathbb{E}_P[\lim_{n \rightarrow \infty} \mathbb{E}_P[\ell(X, Y_n(Z))|Z]] \quad (74)$$

$$= \mathbb{E}_P[-V(P_{X|Z})] \quad (75)$$

$$= -\sum_z P_Z(z)V(P_{X|Z}). \quad (76)$$

Here in (74) we have used the monotone convergence theorem.

The arguments above show that there exists a sequence of measurable functions of Z , $\{Y_n(Z)\}$ such that

$$\lim_{n \rightarrow \infty} \mathbb{E}_P[\ell(X, Y_n(Z))] = -\sum_z P_Z(z)V(P_{X|Z}). \quad (77)$$

Combining it with (70), by the definition of infimum, we know that (65) holds. The claim follows from plugging (63) and (65) into the definition of $C(\ell, P_{XZ})$.

B. Proof of Lemma 2

By Lemma 1, we know there exists a convex function $V : \Gamma_n \rightarrow \mathbb{R}$, such that

$$C(\ell, P_{XZ}) = \left(\sum_z P_Z(z) V(P_{X|Z}) \right) - V(P_X). \quad (78)$$

Let $\delta_i \triangleq (0, 0, \dots, 1, \dots, 0)$ be a distribution in Γ_n that puts mass one on the i -th symbol of \mathcal{X} . Define $a_i \triangleq V(\delta_i)$. We know that $a_i \in \mathbb{R}, \forall i = 1, 2, \dots, n$.

Define the convex function $G : \Gamma_n \rightarrow \mathbb{R}$ as

$$G(P) = V(P) - \sum_{i=1}^n a_i p_i. \quad (79)$$

Now it is easy to verify that $G(\delta_i) = 0, \forall i = 1, 2, \dots, n$. After some algebra we can show that

$$C(\ell, P_{XZ}) = \left(\sum_z P_Z(z) G(P_{X|Z}(\cdot|z)) \right) - G(P_X). \quad (80)$$

Taking $Z \in \mathcal{X}$, and $P_Z = (p_1, p_2, \dots, p_n)$ to be an arbitrary probability distribution. Setting $P_{X|Z}(\cdot|z) = \delta_z$, then we have

$$C(\ell, P_{XZ}) = -G(P_X) = -G((p_1, p_2, \dots, p_n)). \quad (81)$$

Define $T = \pi(X)$ to be a permutation of X , which is sufficient for Z . The Data Processing Axiom implies that

$$C(\ell, P_{XZ}) = C(\ell, P_{TZ}), \quad (82)$$

By construction, we have

$$C(\ell, P_{XZ}) = -G((p_1, p_2, \dots, p_n)), \quad (83)$$

$$C(\ell, P_{TZ}) = -G((p_{\pi^{-1}(1)}, p_{\pi^{-1}(2)}, \dots, p_{\pi^{-1}(n)})), \quad (84)$$

which implies that the function G is invariant to permutations. In other words, G is a symmetric function.

We take $c = -(a_1, a_2, \dots, a_n)$ to finish the proof.